

The Goodness of Covariance Selection Problem from AUC Bounds

Navid Tafaghodi Khajavi and Anthony Kuh

Department of Electrical Engineering,
University of Hawaii, Honolulu, HI 96822
Email: navidt@hawaii.edu, kuh@hawaii.edu

Abstract—We conduct a study of graphical models and discuss the quality of model selection approximation by formulating the problem as a detection problem and examining the area under the curve (AUC). We are specifically looking at the model selection problem for jointly Gaussian random vectors. For Gaussian random vectors, this problem simplifies to the covariance selection problem which is widely discussed in literature by Dempster [1]. In this paper, we give the definition for the correlation approximation matrix (CAM) which contains all information about the model selection problem and discuss the p th order Markov chain model and the p th order star network model for the a Gaussian distribution with Toeplitz covariance matrix. For each model, we compute the model covariance matrix as well as the KL divergence between the Gaussian distribution and its model. We also show that if the model order, p , is proportional to the number of nodes, n , then the model selection is asymptotically good as the number of nodes, n , goes to infinity since the AUC in this case is bounded away from one. We conduct some simulations which confirm the theoretical analysis and also show that the selected model quality increases as the model order, p , increases.

I. INTRODUCTION

In signal processing and machine learning a fundamental problem is to balance performance quality (i.e. minimizing cost function) with computational complexity. A powerful tool in order to address this trade-off is graphical model selection. Model selection methods provide approximated models with desired accuracy as needed for different applications. Given data, different model selection algorithms impose different structure to model data. In the case of jointly Gaussian data, covariance selection problem is presented and studied in [1] and [2]. The purpose of the covariance selection problem is to reduce the computation complexity in various applications.

Some of the model selection algorithms to impose structure are the Chow-Liu minimum spanning tree (MST) [3], the first order Markov chain approximation [4] and penalized likelihood methods such as LASSO [5] and graphical LASSO [6] that can be used to approximate the correlation matrix and inverse correlation matrix with a more sparse graph while retaining good accuracy. The Chow-Liu MST algorithm for Gaussian distribution is to find the optimal tree structure using a Kullback-Leibler (KL) divergence cost function [1]. The Chow-Liu algorithm utilizes the Kruskal algorithm [7]. The first order Markov chain approximation uses a regret cost function to output a chain structured graph [4]. Penalized likelihood methods specify the graph representation by eliminating some of the edges.

In this paper we extend work of [8] where we formulated a covariance model selection paper using a detection problem formulation. The [8] focused on examples where approximation were trees. Here we extend approximations to clique graphs with junction trees. We consider a simple example where the covariance matrix is a Toeplitz covariance matrix with ones along the diagonal and correlation coefficient ρ on the off-diagonals. This covariance matrix is interesting and arise in different applications¹. Given this covariance matrix, we ask the following question, "when is a covariance selection approximation good?" To answer this question we use the detection problem formulation proposed in [8]. The detection problem for Gaussian data leads to calculation of the log-likelihood ratio test (LLRT), the receiver operating characteristic (ROC) curve, the KL divergence and the reverse KL divergence as well as the area under the curve (AUC) where the AUC is used as the accuracy measure for the detection problem on average. We also present the correlation approximation matrix (CAM) as the product of the original correlation matrix and the inverse of the model approximation correlation matrix. For Gaussian data this matrix contains all the information needed to compute the information divergences, the ROC curve and the area under this curve, i.e. the AUC. We present an analytical expression to compute the KL divergence between the original distribution and the model covariance matrix of order, p . We show that if we pick a model order, p , proportional to the number of nodes, n , the AUC is asymptotically bounded away from one as n goes to infinity. Moreover, we present some simulation results. We pick different values as the order of the approximation model and compare the p th order Star approximation model with the p th order Markov chain approximation model. Simulation results show that the p th order star approximation model has smaller AUC than the p th order Markov chain approximation model and thus has better performance. Also, through simulations we confirm our theoretical results showing that the AUC is bounded away from one when model order, p , is proportional to the number of nodes, n .

The rest of this paper is organized as follows. In section II we give the detection problem framework, the sufficient test statistic and the log-likelihood ratio test. Moreover, the

¹Looking at the solar irradiation datasets [9], we can see that sensors that are distributed in small geographical areas are highly correlated and have approximately the same correlations.

sufficient test statistic for Gaussian data as well as its distribution under both hypotheses are also presented in this section. The ROC curve and the definition of AUC as well as analytical expression for the AUC are presented in this section. Section III provides the theoretical analysis of the Toeplitz covariance matrix with ones along the diagonal and correlation coefficient ρ 's on the off-diagonals. The model covariance matrix for a given order, p , as well as the KL divergence between the original distribution and the model distribution are also presented in this section. Moreover, asymptotic upper bounds for KL divergence and the AUC are also presented in this section. In section IV we present some simulation results for approximation model with different orders and investigates the quality of different model approximations based on the numerically evaluated AUC and also its analytical upper and lower bounds. Finally, Section V summarizes results of this paper and discuss further research directions.

Notation remark: In the rest of this paper, with abuse of notion, when we use the KL divergence between random vectors it means the KL divergence between their associated distributions.

II. DETECTION PROBLEM FRAMEWORK

A. Preliminaries

Let $\underline{X} \sim \mathcal{N}(\underline{0}, \Sigma_{\underline{X}})$, i.e. jointly Gaussian with mean 0 and covariance matrix $\Sigma_{\underline{X}}$, where $\underline{X} \in \mathbb{R}^n$. We want to approximate the random vector \underline{X} , with another random vector, $\underline{X}_{\mathcal{M}} \in \mathbb{R}^n$ which has a zero-mean jointly Gaussian distribution with the covariance matrix $\Sigma_{\underline{X}_{\mathcal{M}}}$ associated with the desired model², i.e. $\underline{X} \sim \mathcal{N}(\underline{0}, \Sigma_{\underline{X}_{\mathcal{M}}})$. Note that the model covariance matrix is also positive-definite, $\Sigma_{\underline{X}_{\mathcal{M}}} > 0$. Also, let $\mathcal{G} = (\mathcal{V}, \mathcal{E}_{\mathcal{M}})$ be the graph representation of the model random vector $\underline{X}_{\mathcal{M}}$ where sets \mathcal{V} and $\mathcal{E}_{\mathcal{M}} \subseteq \psi$ are the set of all vertices and the set of all edges of the graph representing of $\underline{X}_{\mathcal{M}}$, respectively where ψ is the set of all edges in complete graph with vertex set \mathcal{V} .

We define the correlation approximation matrix (CAM) associated with the model selection problem as follows.

Definition 1. Correlation approximation matrix [8]. The CAM for the model is defined as $\Delta \triangleq \Sigma_{\underline{X}} \Sigma_{\underline{X}_{\mathcal{M}}}^{-1}$. ■

Remark: The CAM is a positive definite matrix and its eigenvalues contains all information necessary to compute cost functions associated with the model selection problem.

B. General Framework

A common measure to compare two probability distributions is the KL divergence. Here we expand the comparison by considering a detection problem where the null hypothesis represents the original random vector, \underline{X} and the alternate hypothesis represents the approximate random vector $\underline{X}_{\mathcal{M}}$. We need to define a test statistic to quantify the detection problem. The likelihood ratio test (the Neyman-Pearson (NP)

Lemma [10]) is the most powerful test statistic where we first define the log-likelihood ratio test (LLRT) as

$$l(\underline{x}) = \log \frac{f_{\underline{X}}(\underline{x}|\mathcal{H}_1)}{f_{\underline{X}}(\underline{x}|\mathcal{H}_0)}$$

where $f_{\underline{X}}(\underline{x}|\mathcal{H}_0)$ is the distribution of random vector \underline{X} under the null hypothesis while $f_{\underline{X}}(\underline{x}|\mathcal{H}_1)$ is the distribution of random vector \underline{X} under the alternative hypothesis. Moreover, let $L(\underline{X})$ be the LLRT random variable. Also, let random variables

$$L_0 \triangleq L(\underline{X})|\mathcal{H}_0$$

and

$$L_1 \triangleq L(\underline{X})|\mathcal{H}_1$$

be the LLRT statistics under hypothesis \mathcal{H}_0 and hypothesis \mathcal{H}_1 , respectively. We then define the *false-alarm probability* and the *detection probability* by comparing the LLRT statistic under each hypothesis with a given threshold, τ , and computing the following probabilities

- The false-alarm probability, $P_0(\tau)$, under the null hypothesis, \mathcal{H}_0 : $P_0(\tau) = \Pr(L_0 \geq \tau)$,
- The detection probability, $P_1(\tau)$, under the alternative hypothesis, \mathcal{H}_1 : $P_1(\tau) = \Pr(L_1 \geq \tau)$.

The most powerful test is defined by setting the false-alarm rate $P_0(\tau) = \bar{P}_0$ and then computing the threshold value τ_0 such that $\Pr(L_0 \geq \tau_0) = \bar{P}_0$.

Definition 2. The KL divergence between two multivariate continuous distributions $p(\underline{X})$ and $q(\underline{X})$ is defined as

$$\mathcal{D}(p_{\underline{X}}(\underline{x})||q_{\underline{X}}(\underline{x})) = \int_{\mathcal{X}} p_{\underline{X}}(\underline{x}) \log \frac{p_{\underline{X}}(\underline{x})}{q_{\underline{X}}(\underline{x})} d\underline{x}$$

where \mathcal{X} is the feasible set. ■

Throughout this paper we may use other notations such as the KL divergence between two random vectors or the KL divergence between two covariance matrices for zero-mean Gaussian distribution case in order to present the KL divergence between two distributions.

Proposition 1. Expectation of the LLRT statistic under each hypothesis is

- $E(L_0) = -\mathcal{D}(f_{\underline{X}}(\underline{x}|\mathcal{H}_0)||f_{\underline{X}}(\underline{x}|\mathcal{H}_1))$,
- $E(L_1) = \mathcal{D}(f_{\underline{X}}(\underline{x}|\mathcal{H}_1)||f_{\underline{X}}(\underline{x}|\mathcal{H}_0))$.

Proof: Proof is based on the KL divergence definition. ■

The NP decision rule in a regular detection problem framework is to accept the hypothesis \mathcal{H}_1 if the LLRT statistic, $L(\underline{x})$, exceeds a critical value which is set based on the false-alarm probability, and reject it otherwise. As it is mentioned in [8], we pursue a different goal in the approximation problem scenario. We approximate a model distribution, $f_{\underline{X}_{\mathcal{M}}}(\underline{x})$, as close as possible to the given distribution, $f_{\underline{X}}(\underline{x})$. In ideal case where there is no approximation error, the detection probability must be equal to the false-alarm probability for the optimal detector at all possible thresholds, i.e. the receiver operating characteristic (ROC) curve [11] that represents best

²Examples of possible models: star structure and Markov chain.

detectors for all threshold values should be a line of slope 1 passing through the origin.

C. Multivariate Gaussian distribution

Let the random vector \underline{X} have a multivariate Gaussian distribution with covariance matrix $\Sigma_{\underline{X}}$. In this paper, the null hypothesis, \mathcal{H}_0 , is the hypothesis that the parameter of interest, which is the covariance matrix of the random vector \underline{X} , is known and is equal to $\Sigma_{\underline{X}}$ while the alternative hypothesis, \mathcal{H}_1 , is the hypothesis that the random vector \underline{X} is replaced by the model random vector $\underline{X}_{\mathcal{M}}$ which means that the random vector \underline{X} has the model approximation distribution with the covariance matrix, $\Sigma_{\underline{X}_{\mathcal{M}}}$. Thus, we can rewrite the LLRT statistic as

$$l(\underline{x}) = \log \frac{f_{\underline{X}_{\mathcal{M}}}(\underline{x})}{f_{\underline{X}}(\underline{x})}$$

The LLRT statistic can be simplified for the multivariate Gaussian distributed random vectors as

$$l(\underline{x}) = \log \frac{\mathcal{N}(\underline{0}, \Sigma_{\underline{X}_{\mathcal{M}}})}{\mathcal{N}(\underline{0}, \Sigma_{\underline{X}})} = -c + k(\underline{x}) \quad (1)$$

where $c = -\frac{1}{2} \log(|\Delta|)$ is a constant and $k(\underline{x}) = \underline{x}^T \mathbf{K} \underline{x}$ where $\mathbf{K} = \frac{1}{2}(\Sigma_{\underline{X}}^{-1} - \Sigma_{\underline{X}_{\mathcal{M}}}^{-1})$ is an indefinite matrix with both positive and negative eigenvalues.

Theorem 1. Covariance Selection [1]. *Given a multivariate Gaussian distribution with covariance matrix $\Sigma_{\underline{X}} > 0$, $f_{\underline{X}}(\underline{x})$, and a model \mathcal{M} , there exists a unique approximated multivariate Gaussian distribution with covariance matrix $\Sigma_{\underline{X}_{\mathcal{M}}} > 0$, $f_{\underline{X}_{\mathcal{M}}}(\underline{x})$, that minimize the KL divergence, $\mathcal{D}(f_{\underline{X}}(\underline{x}) || f_{\underline{X}_{\mathcal{M}}}(\underline{x}))$ and satisfies the covariance selection rules, i.e. the model covariance matrix satisfies the following covariance selection rules*

- $\Sigma_{\underline{X}_{\mathcal{M}}}(i, i) = \Sigma_{\underline{X}}(i, i), \quad \forall i \in \mathcal{V}$
- $\Sigma_{\underline{X}_{\mathcal{M}}}(i, j) = \Sigma_{\underline{X}}(i, j), \quad \forall (i, j) \in \mathcal{E}_{\mathcal{M}}$
- $\Sigma_{\underline{X}_{\mathcal{M}}}^{-1}(i, j) = 0, \quad \forall (i, j) \in \mathcal{E}_{\mathcal{M}}^c$

where the set $\mathcal{E}_{\mathcal{M}}^c = \psi - \mathcal{E}_{\mathcal{M}}$ represents the complement of the set $\mathcal{E}_{\mathcal{M}}$.

Proof: Proof for Gaussian distributions is given in Dempster 1972 paper [1]. ■

Remark: From theorem 1 and definition of the KL divergence for Gaussian distributions, we have $c = \mathcal{D}(f_{\underline{X}}(\underline{x}) || f_{\underline{X}_{\mathcal{M}}}(\underline{x}))$, since given any covariance matrix and its model covariance matrix satisfying theorem 1, we have $\text{tr}(\Delta) = n$.

D. Distribution of the LLRT statistic

The random vector \underline{X} has Gaussian distribution under both hypotheses \mathcal{H}_0 and \mathcal{H}_1 . Thus under both hypotheses, the real random variable, $K(\underline{X}) \triangleq \underline{X}^T \mathbf{K} \underline{X}$ has a generalized chi-squared distribution, i.e. the random variable, $K(\underline{X})$, is equal to a weighted sum of chi-squared random variables with both positive and negative weights under both hypotheses. Let us define $\underline{W} = \Sigma_{\underline{X}}^{-\frac{1}{2}} \underline{X}$ under \mathcal{H}_0 and $\underline{Z} = \Sigma_{\underline{X}_{\mathcal{M}}}^{-\frac{1}{2}} \underline{X}$ under \mathcal{H}_1 , where $\Sigma_{\underline{X}}^{\frac{1}{2}}$ and $\Sigma_{\underline{X}_{\mathcal{M}}}^{\frac{1}{2}}$ are the square root of covariance

matrices $\Sigma_{\underline{X}}$ and $\Sigma_{\underline{X}_{\mathcal{M}}}$, respectively. Then let random vectors $\underline{W} \sim \mathcal{N}(\underline{0}, \mathbf{I})$ and $\underline{Z} \sim \mathcal{N}(\underline{0}, \mathbf{I})$ have zero-mean Gaussian distributions with the same covariance matrices, \mathbf{I} , where \mathbf{I} is the identity matrix of dimension n . Note that, the CAM is a positive definite matrix with $\lambda_i > 0$ where $1 \leq i \leq n$. Thus, the random variable $K(\underline{X})$, under both hypotheses \mathcal{H}_0 and \mathcal{H}_1 can be defined as

$$K_0 \triangleq \frac{1}{2} \sum_{i=1}^n (1 - \lambda_i) W_i^2$$

and

$$K_1 \triangleq \frac{1}{2} \sum_{i=1}^n (\lambda_i^{-1} - 1) Z_i^2$$

respectively, where random variables W_i and Z_i , are the i -th element of random vectors \underline{W} and \underline{Z} , respectively. Moreover, random variables W_i^2 and Z_i^2 , follow the first order central chi-squared distribution. Note that, $L_0 = -c + K_0$ and $L_1 = -c + K_1$.

Remark: As a simple consequence of the covariance selection theorem, the summation of weights for the generalized chi-squared random variable, the expectation of $K(\underline{X})$, is zero under the hypothesis \mathcal{H}_0 , i.e. $E(K_0) = \frac{1}{2} \sum_{i=1}^n (1 - \lambda_i) = 0$ [1], and this summation is positive under the hypothesis \mathcal{H}_1 , i.e. $E(K_1) = \frac{1}{2} \sum_{i=1}^n (\lambda_i^{-1} - 1) \geq 0$.

E. Area under the curve

As we mentioned before, in approximation set up the desired goal is that the ROC curve is as close as possible to the line of slope 1 passing through the origin in comparison to the step function in the hypothesis testing problem [8]. Area under the curve is defined as the integral of the ROC curve. Note that in approximation problem presented here we want it to be around half. Area under the curve (AUC) is defined as the integral of the ROC curve (figure 1) and is a measure of accuracy in decision problems.

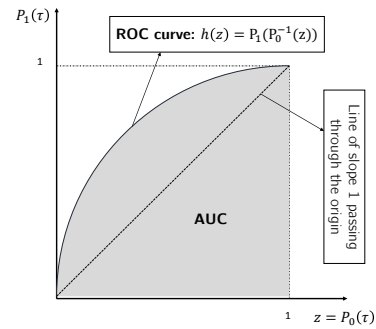


Fig. 1: The ROC curve and the area under the ROC curve. Each point on the ROC curve indicates a detector with given detection and false-alarm probabilities.

Definition 3. *The area under the ROC curve (AUC) is defined as*

$$AUC = \int_0^1 h(z) dz = \int_0^1 P_1(\tau) dP_0(\tau) \quad (2)$$

where τ is the detection problem threshold. ■

Remark: The AUC is a measure of accuracy for the detection problem and $1/2 \leq \text{AUC} \leq 1$. Note that, in conventional decision problems, the AUC is desired to be as close as possible to 1 while in approximation problem presented here we want the AUC to be close to $1/2$.

Theorem 2. Statistical property of AUC [12]. *The AUC for the LLRT statistic, $L(\underline{X})$, and two hypotheses, \mathcal{H}_0 and \mathcal{H}_1 is*

$$\text{AUC} = \Pr(L_\Delta > 0).$$

where $L_\Delta \triangleq L_1 - L_0$. ■

III. TOEPLITZ COVARIANCE MATRIX

Here, we assumed that the n by n covariance matrix $\Sigma_{\underline{X}}$ has a Toeplitz structure with ones on the diagonal and the correlation coefficient ρ as off diagonal elements

$$\Sigma_{\underline{X}} = \begin{bmatrix} 1 & \rho & \dots & \rho \\ \rho & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \rho \\ \rho & \dots & \rho & 1 \end{bmatrix}.$$

Definition 4. Clique. *A maximal subset of the nodes which defines a complete subgraph is the clique subgraph.* ■

In other words, all pairs of nodes are connected in the clique subgraph.

Definition 5. Junction tree. *A junction tree is a clique tree [13] such that for each pair of cliques C_1 and C_2 in the graph, all cliques on the path between C_1 and C_2 contain their intersection, $C_1 \cap C_2$.* ■

In this example, we are interested in models which can be represented using junction trees whose vertices are cliques of the size at most p .³ Going back to the model selection problem for the example, we are investigating the following two generalizations of the chain and the star networks. Note that, we can construct a junction tree for these two special models.

A. p th order star network

The model covariance matrix for the p th order star network where all nodes are connected to the first p nodes which all are connected together is as follow

$$\Sigma_{\underline{X}_{\mathcal{M}}}^{pth-star} = \begin{bmatrix} 1 & \rho & \dots & \dots & \dots & \dots & \rho \\ \rho & \ddots & \ddots & & & & \vdots \\ \vdots & \ddots & 1 & \rho & \dots & \dots & \rho \\ \vdots & & \rho & 1 & \rho_1 & \dots & \rho_1 \\ \vdots & & \vdots & \rho_1 & \ddots & \ddots & \vdots \\ \vdots & & \vdots & \vdots & \ddots & \ddots & \rho_1 \\ \rho & \dots & \rho & \rho_1 & \dots & \rho_1 & 1 \end{bmatrix}$$

³We avoid cycles by turning subsets of the nodes into supernodes.

where

$$\rho_1 = \frac{p\rho^2}{(p-1)\rho + 1}.$$

B. p th order Markov chain network

The model covariance matrix for the p th order Markov chain network is as follow

$$\Sigma_{\underline{X}_{\mathcal{M}}}^{pth-chain} = \begin{bmatrix} 1 & \rho & \dots & \rho & \rho_1 & \dots & \rho_{n-p-1} \\ \rho & \ddots & \ddots & & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & & \ddots & \rho_1 \\ \rho & & \ddots & \ddots & \ddots & & \rho \\ \rho_1 & & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & \ddots & \rho \\ \rho_{n-p-1} & \dots & \rho_1 & \rho & \dots & \rho & 1 \end{bmatrix}.$$

To satisfy Theorem 1 we have that ρ_i for $i \in \{1, \dots, n-p-1\}$ can be computed through the following recursive equation

$$\rho_i = \rho_{i-1}^T \underline{v}_i \frac{\rho}{(p-1)\rho + 1} \quad (3)$$

where $\underline{v}_i = [\overbrace{1, \dots, 1}^p, 0, \dots, 0]^T$ is a vector of length n and $\underline{\rho}_i = [\rho_i, \dots, \rho_1, \overbrace{\rho, \dots, \rho}^p]^T$ where $\underline{\rho}_0 = [\overbrace{\rho, \dots, \rho}^p]^T$ is the initialization step.

Lemma 3. *The KL divergence for the p th order star network and the p th order Markov chain network can be calculated as*

$$\begin{aligned} \mathcal{D}(\underline{X} || \underline{X}_{pth-chain}) &= \frac{1}{2}(n-p) \log \left(\frac{p\rho + 1}{(p-1)\rho + 1} \right) \\ &\quad + \frac{1}{2} \log \left(\frac{(p-1)\rho + 1}{(n-1)\rho + 1} \right) \end{aligned}$$

and

$$\mathcal{D}(\underline{X} || \underline{X}_{pth-star}) = \mathcal{D}(\underline{X} || \underline{X}_{pth-chain}).$$

Proof: Note that, from [14] we have

$$|\Sigma_{\underline{X}_{\mathcal{M}}}^{pth-chain}| = \frac{[(p\rho + 1)(\rho - 1)^p]^{(n-p)}}{[((p-1)\rho + 1)(\rho - 1)^{p-1}]^{(n-p-1)}}$$

and

$$|\Sigma_{\underline{X}}| = ((n-1)\rho + 1)(\rho - 1)^{n-1}.$$

Inserting the values of these determinants into the KL divergence

$$\mathcal{D}(\underline{X} || \underline{X}_{\mathcal{M}}) = -\frac{1}{2} \log \left(\Sigma_{\underline{X}} \Sigma_{\underline{X}_{\mathcal{M}}}^{-1} \right)$$

we conclude the result for the p th order Markov chain network. To show that the KL divergence for the p th order star network is exactly equal to the KL divergence for the p th order chain network, we need to construct the corresponding junction tree for each of these networks by grouping appropriate p nodes. Note that, the KL divergence for the junction trees are equal since the mutual information between the junction nodes are exactly equal. ■

Proposition 2. The KL divergence for the p th order star network and the p th order Markov chain network is bounded as n goes to infinity if for a given constant number, $\kappa > 1$, the order, p , is the integer number in interval,

$$\mathcal{D}(\underline{X}||\underline{X}_{pth-star}) < \infty \quad \text{as} \quad (n \rightarrow \infty, n/p \rightarrow \kappa).$$

Proof: Let $p = \lceil n/\kappa \rceil$ be the smallest integer greater than or equal to n/κ . The KL divergence can be bounded as follow

$$\begin{aligned} \mathcal{D}(\underline{X}||\underline{X}_{pth-star}) &= \frac{(n - \lceil n/\kappa \rceil)}{2} \log \left(1 + \frac{\rho}{(\lceil n/\kappa \rceil - 1)\rho + 1} \right) \\ &\quad + \frac{1}{2} \log \left(\frac{(\lceil n/\kappa \rceil - 1)\rho + 1}{(n - 1)\rho + 1} \right) \\ &\stackrel{(a)}{\leq} \frac{(n - n/\kappa)}{2} \log \left(1 + \frac{\rho}{(n/\kappa - 1)\rho + 1} \right) \\ &\quad + \frac{1}{2} \log \left(\frac{((n/\kappa + 1) - 1)\rho + 1}{(n - 1)\rho + 1} \right) \\ &\stackrel{(b)}{\leq} \frac{(1 - 1/\kappa)n}{2} \left(\frac{\rho}{(n/\kappa - 1)\rho + 1} \right) \\ &\quad + \frac{1}{2} \log \left(\frac{(n/\kappa)\rho + 1}{(n - 1)\rho + 1} \right) \end{aligned}$$

Where (a) is true since for the integer order, p , we have $n/\kappa \leq p < n/\kappa + 1$ and (b) is true since $\log(1 + z) \leq z$ for $z \geq 0$. Then, in the limit we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathcal{D}(\underline{X}||\underline{X}_{pth-star}) &\leq \frac{(1 - 1/\kappa)}{2/\kappa} + \frac{1}{2} \log(1/\kappa) \\ &\leq \frac{\kappa - 1}{2} - \frac{\log(\kappa)}{2} < \infty \end{aligned}$$

which complete the proof. ■

Proposition 3. The AUC of the p th order star network and the p th order Markov chain network is bounded from 1 as n goes to infinity if $n = \kappa p$ and $p = \lceil n/\kappa \rceil$,

$$\Pr(L_{\Delta} > 0) < 1.$$

Proof: We can conclude this result from the proposition 2 upper bound for the KL divergence combined with the upper bound for the AUC,

$$\Pr(L_{\Delta} > 0) \leq 1 - e^{-\lim_{n \rightarrow \infty} \mathcal{D}(\underline{X}||\underline{X}_{pth-star}) - 1} < 1$$

provided in [8]. ■

IV. SIMULATION RESULTS AND DISCUSSION

In this section, we consider the Toeplitz example presented before as the covariance matrix for a Gaussian random vector. We calculate different models such as the p th order Markov chain and the p th order star networks for various values of p . For a given order, both of the aforementioned models have the same KL divergence values as calculated in lemma 3. Moreover, we compute AUC and compare it with its lower and upper bounds [8] for these cases.

Figure 2 plots $(1 - \text{AUC})$ in log-scale v.s. the dimension of the graph, n , in linear-scale for star approximation (**left**) and chain approximation (**right**) with different model orders, $p = 1, p = 3, p = 5$ and $p = 7$ for correlation coefficient $\rho = 0.9$. As it is indicated in this figure, $(1 - \text{AUC})$ decreases as the order of the model increases for both star and chain models. Moreover, from this figure, we can conclude that the p th order star network performs better than the p th order Markov chain network since $(1 - \text{AUC})$ decay exponent is smaller for the former model than the latter model. This can also be seen by comparing the covariance matrix $\Sigma_{\underline{X}}$ and the model covariance matrix, $\Sigma_{\underline{X}_{\mathcal{M}}}$ where the model covariance matrix associated with the p th order star network is more similar to the covariance matrix $\Sigma_{\underline{X}}$ than the model covariance matrix associated with the p th order Markov chain network. For example, even the quality of the first order star network approximation is better than the quality of the fifth order Markov chain approximation in the simulation results provided in this figure.

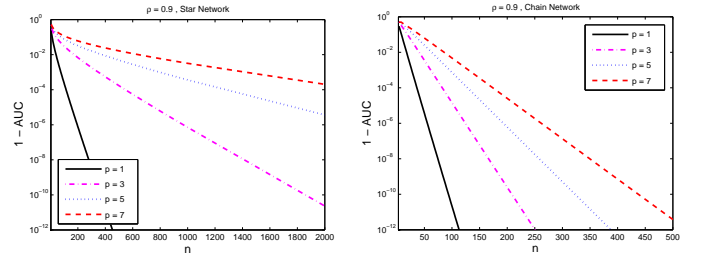


Fig. 2: $1 - \text{AUC}$ (log-scale) v.s. the dimension of the graph (linear-scale), n , for star approximation (**left**) and chain approximation (**right**) with different model orders, $p = 1, p = 3, p = 5$ and $p = 7$ and correlation coefficient $\rho = 0.9$.

Figure 3 plots KL divergence v.s. $-\log(1 - \text{AUC})$ for the presented models. In this figure, the dimension n is set to 15, the order p is set to 1 and 3 and the correlation coefficient ρ is set to 0.9. Furthermore, the feasible region presented in [8] and its asymptotic behavior are also plotted in this figure. For both models, the KL divergence and the reverse KL divergence are computed and are plotted on this figure. Note that, KL divergences for both models are equal (see lemma 3) and are connected in this figure. As it is shown in the figure, the third order model has better performance than the first order model.

Figure 4 plots $1 - \text{AUC}$ v.s. the dimension of the graph, n for the p th order star approximation of the Toeplitz example for $\rho = 0.1$ (**left**) and $\rho = 0.9$ (**right**) while keeping the model order proportional to the number of nodes in the graphical model, n . More specifically, in this figure, we set the model order $p = \lceil n/\kappa \rceil$ where $\kappa = 10$. Moreover, this figure plots the lower bound and the upper bound for $1 - \text{AUC}$ ⁴. From this figure, we conclude that, p th order star approximation is a good approximation model when the model order, p is proportional to the number of nodes, n , since the AUC is bounded from

⁴Bounds are presented in [8].

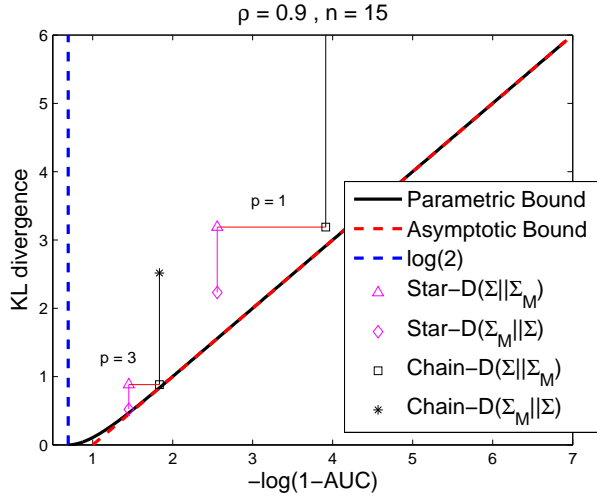


Fig. 3: KL divergence v.s. AUC and the AUC parametric bound [8] v.s. for graph dimension, $n = 15$ for the p th order Markov chain approximation and p th order star network for $p = 1$ and $p = 3$ with $\rho = 0.9$.

one as $n \rightarrow \infty$. Similarly, figure 5 plots $1 - \text{AUC}$ and its upper and lower bounds v.s. the dimension of the graph, n for the p th order Markov chain approximation of the Toeplitz example for $\rho = 0.1$ (left) and $\rho = 0.9$ (right) with $p = \lceil n/\kappa \rceil$ where $\kappa = 10$. Plots in this figure are not monotonically decreasing since both the order p and the dimension n are integers and thus the ratio n/p is not exactly equal to κ for all values of p and n . Furthermore, from the figure, the p th order Markov chain approximation is a good approximation model when the model order, p is proportional to the number of nodes, n , since the AUC is bounded from one as $n \rightarrow \infty$. Comparing the plots in figure 4 and figure 5 we can clearly see that even though the AUC for both approximation models are bounded from one, the p th order star approximation model is a better model than the p th order Markov chain approximation model.

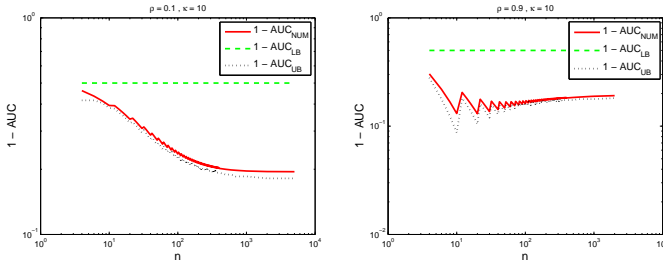


Fig. 4: $1 - \text{AUC}$ and its lower and upper bounds v.s. the dimension of the graph, n for the p th order star approximation of the Toeplitz example for $\rho = 0.1$ (left) and $\rho = 0.9$ (right) with the model order $p = \lceil n/\kappa \rceil$ where $\kappa = 10$.

V. CONCLUSION

In this paper, we formulate a detection problem to investigate the quality of the graphical model approximation.

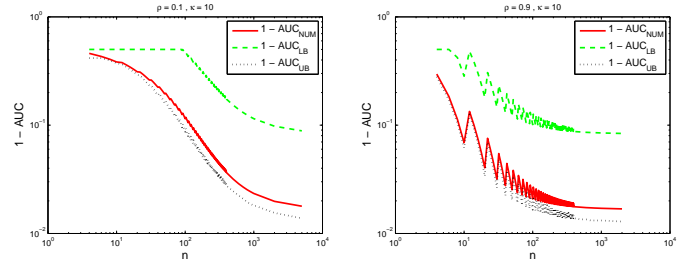


Fig. 5: $1 - \text{AUC}$ and its lower and upper bounds v.s. the dimension of the graph, n for the p th order Markov chain approximation of the Toeplitz example for $\rho = 0.1$ (left) and $\rho = 0.9$ (right) with the model order $p = \lceil n/\kappa \rceil$ where $\kappa = 10$.

We discuss the quality of model selection approximation by examining the area under the curve (AUC). We consider jointly Gaussian random vectors and give the definition for the correlation approximation matrix (CAM). We discuss graphical models with junction trees such as the p th order Markov chain and the corresponding star network interpretation for a special Toeplitz covariance matrix with ones along the diagonal and correlation coefficient ρ 's on the off-diagonals. These models has very short loops and has associated junction tree that connects cliques of the same size. The model covariance matrix as well as the KL divergence between the original distribution and the model distribution are computed for the presented Toeplitz covariance matrix. We also quantify the goodness of the covariance selection problem for this Toeplitz covariance matrix. For this covariance matrix, we show that if the model order, p , is proportional to the number of nodes, n , then the model selection is asymptotically good as $n \rightarrow \infty$ since the AUC is asymptotically bounded away from one. We conduct some simulations which show that the selected model quality increases as the model order, p , increases and confirm our theoretical results.

ACKNOWLEDGMENT

This work was supported in part by NSF grant ECCS-1310634, and the University of Hawaii REIS project.

REFERENCES

- [1] A. P. Dempster, "Covariance selection," *Biometrics*, vol. 28, no. 1, pp. 157–175, March 1972.
- [2] Steffen L Lauritzen, *Graphical models*, Clarendon Press, 1996.
- [3] C. K. Chow and C. N. Liu, "Approximating discrete probability distributions with dependence trees," *IEEE Transactions on Information Theory*, pp. 462–467, 1968.
- [4] N. T. Khajavi and A. Kuh, "First order markov chain approximation of microgrid renewable generators covariance matrix," in *Proc. of IEEE International Symposium on Information Theory, Istanbul, Turkey (ISIT'13)*, July 2013, pp. 1207–1211.
- [5] N. Meinshausen and P. Bühlmann, "Model selection through sparse maximum likelihood estimation," *Annals of Statistics*, pp. 1436–1464, 2006.
- [6] Jerome Friedman, Trevor Hastie, and Robert Tibshirani, "Sparse inverse covariance estimation with the graphical lasso," *Biostatistics*, vol. 9, no. 3, pp. 432–441, 2008.
- [7] J. B. Kruskal, "On the shortest spanning subtree of a graph and the traveling salesman problem," *Proceedings of the American Mathematical society*, vol. 7, no. 1, pp. 48–50, 1956.

- [8] Navid Tafaghodi Khajavi and Anthony Kuh, "The quality of the covariance selection through detection problem and auc bounds," *arXiv preprint arXiv:1605.05776*, 2016.
- [9] N. T. Khajavi, A. Kuh, and N. P. Santhanam, "Spatial correlations for solar pv generation and its tree approximation analysis," in *Proc. of the Asia-Pacific Signal and Information Processing Association (APSIPA ASC)*, Dec 2014, pp. 1–5.
- [10] J. Neyman and E. S. Pearson, "On the use and interpretation of certain test criteria for purposes of statistical inference," *Biometrika*, vol. 20, 1928.
- [11] L. L. Scharf, *Statistical signal processing*, vol. 98, Addison-Wesley Reading, MA, 1991.
- [12] J. A. Hanley and B. J. McNeil, "The meaning and use of the area under a receiver operating characteristic (roc) curve.," *Radiology*, vol. 143, no. 1, pp. 29–36, 1982.
- [13] Jean RS Blair and Barry Peyton, "An introduction to chordal graphs and clique trees," in *Graph theory and sparse matrix computation*, pp. 1–29. Springer, 1993.
- [14] A. Kavcic and J. M. F. Moura, "Matrices with banded inverses: Inversion algorithms and factorization of gauss-markov processes," *IEEE Transactions on Information Theory*, vol. 46, pp. 1495–1509, July 2000.